

Identifying Complex Networks by Random Walks

Filipi Nascimento Silva and Luciano da Fontoura Costa *

December 14, 2006

Abstract

The possibility to identify the nature (e.g. random or scale free) of complex networks while performing respective random walks is investigated with respect to autonomous agents based on Bayesian decision theory and humans navigating through a graphic-interactive interface. The results indicate that the type of the network (choice between random and scale free models) can be correctly estimated in most cases.

‘They rebuild Ersilia elsewhere. They weave a similar pattern of strings which they would like to be more complex and at the same time more regular than the other. (I. Calvino, Invisible Cities’

1 Introduction

Conscious human existence can be understood as a trajectory in the space-time phase space, unfolding as a consequence of our perceptions, decisions and actions. As suggested recently [1, 2], the dynamic evolution of the life of a human individual can be approximated by a random walk¹ [1, 2, 3, 4, 5, 6] in a complex network Γ (e.g. [7, 8, 9]) where the nodes correspond to the possible decisions and the links to the transitions between such decisions. Note that, in case time is taken explicitly, such a complex network will correspond to a decision tree and include no cycles (i.e. a closed path). In order to allow the formation of cycles, we henceforth consider the evolution of time implicitly, allowing that oneself will find her/himself in recurring situations (e.g. choosing shoes for dinner). A series of interesting insights and results about our perception of the complexity of our individual life (as far as decisions are concerned) can be achieved by considering such a model. For instance, in case the complex network Γ is scale free, the average degree of the nodes

as sampled by a traditional random walk will tend to result twice as large as its real value (e.g. [2]). This interesting phenomenon is a direct consequence of the fact that hubs are more likely to be visited during a random walk in a BA network, hence the overestimation of the average degree.

An important issue related to modeling human experience in terms of random walks in complex networks concerns our ability to identify the most likely mathematical model (e.g. Erdős-Rényi or Barabási-Albert) of the complex network being explored while navigating on it through random walks. The current article explores this key issue from the perspective of having human subjects to navigate through ER and BA complex networks models while trying to identify their type. In addition, in order to gain deeper insight on this problem, the subjects correct ratio is compared with results obtained by a agent that uses a optimal identification algorithm (Bayesian decision).

The article starts by presenting basic concepts in complex networks and follows by summarizing some statistical methods (Pearson correlation coefficient and Bayesian decision theory) and explaining the basic concepts of agent classification of networks. The article follows by describing the experimental methodology and presenting its results and comparison with the autonomous agents results.

2 Networks Generation and Degree Distributions

Consider a matrix K , with $K(i, j) = K(j, i) = \{1 \text{ or } 0\}$ (i.e. a binary and symmetric matrix). A non-weighted and non-oriented network (i.e. a graph) can be completely specified by such an *adjacency matrix*, where the existence of a connection between two nodes (i and j) is represented as $K(i, j) = K(j, i) = 1$. The adjacency matrix of a random network (i.e. ER model) of size $N \times N$ can be generated by starting with all elements equal to zero and making $K(i, j) = K(j, i) = 1$ with probability ρ for every pair of nodes i and j , implying average degree equal to $\langle k \rangle = \rho N$. The average degree distribution obtained for 2000 realizations of the ER model with $N = 100$ and $\langle k \rangle = 4$ is shown in Figure 1.

*Cybernetic Vision Research Group, GII-IFSC, Universidade de São Paulo, São Carlos, SP, Caixa Postal 369, 13560-970, Brasil, luciano@ifsc.usp.br.

¹The term *random* is here meant to express general probabilistic decision models, not necessarily the traditional random walk where decisions are taken by uniformly sampling among the paths emanating from each node.

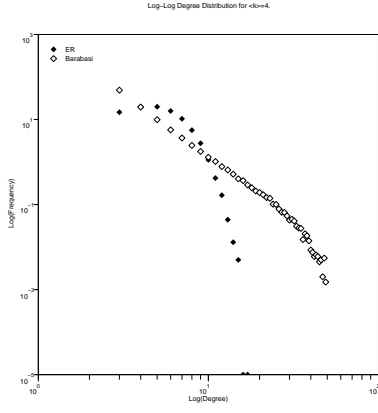


Figure 1: Average degree distributions for 2000 realizations of ER and BA networks with average degree $\langle k \rangle = 4$.

BA networks can be generated by selecting M_0 rows or columns, $K(i, \dots)$ (initial nodes), and then connecting e_0 of these nodes, randomly chosen, with a new node by filling the corresponding value on adjacent matrix [7]. This process is repeated t times, always connecting e_0 of previous selected nodes with a new node.

Therefore, the final number of nodes will be $N = t + M_0$ and the number of edges $e = e_0 t$, with average node degree given as $\langle k \rangle = \frac{2e_0 t}{(M_0 + t)}$. For values of $M_0 \ll t$, the average degree will be $\langle k \rangle = 2e_0$. Note that for low values of M_0 the BA models can only be generated with even average degree.

The degree distribution of BA models is known to follow a power-law [7, 8, 9] as a consequence of its scale-free nature, with power coefficient $\alpha = -3$. The average degree distribution obtained for 2000 realizations of the BA model with $N = 100$ and $\langle k \rangle = 4$ is shown in Figure 1.

3 Statistical Concepts

Two basic statistical methodologies are used in the present work in order to construct an agent that classify a network. Because of the linear behavior of the distribution of average degree for BA model and non-linear for ER model, this property can be exploited as a parameter for the segregation between the two models, this can be made by using the Pearson correlation coefficient as follows.

The Pearson correlation coefficient is a statistical measurement quantifying how strong is the linear joint variation between two random variables [10]. Given the normalized distribution of two random discrete variables, X and Y , the Pearson correlation coefficient

between them is defined by the covariance of that two variables divided by their standard deviation, i.e.:

$$r_{XY} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} \quad (1)$$

Given n samples of the random variables X and Y , henceforth expressed x_i and y_i , the respective Pearson Correlation Coefficient can be estimated as:

$$r_{XY} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{(n - 1)\sigma_X \sigma_Y} \quad (2)$$

Where \bar{x} and \bar{y} are the average values of elements x_i and y_i , and σ_x and σ_y are the respective standard deviations. These conditions bond the Pearson correlation coefficient between -1 and 1. Values of r_{XY} near 0 suggest absence of linear correlation between the two variables, while values around 1 and -1 indicate correlated or anti-correlated behavior, respectively. Note that a nearly straight distribution of points is observed for the cases characterized by absolute values of Pearson coefficients nearly equal to 1.

Because of the linear behavior of the logarithmic node degree distribution observed for scale-free networks, contrasting with the binomial distribution for ER, their squared correlation coefficients can be used as a sound criterion for discrimination and identification of these two models. Therefore, the Pearson correlation coefficient is used in this work in order to quantify the degree of straightness of the log-log degree distribution. As shown in figure 2, the squared correlation coefficients for the logarithmic of node degree distributions for BA model tend to be substantially higher than those for the ER model. Note that the difference between the Pearson coefficients for the BA model decreases as the average degrees of those networks increase, as a consequence of the small size of the adopted networks (i.e. $N=100$).

Bayes decision theory is the optimal statistical method for supervised classification of data, provided the density distribution of the characteristics of the data classes is known (It can be formally verified [12] that the Bayesian decision criterion adopted in this work is optimal from the perspective of minimizing the probability of misclassifications). In the specific case of equiprobable classes, the Bayes decision involves selecting the class that is most probable for a set of measured properties of an element. Suppose we have two equiprobable classes of elements, A and B , and let e be an unknown element whose class must be determined by using some measured property h_e . Provided the density distribution functions, $\rho_A(h)$ and $\rho_B(h)$ are available, Bayes decision theory selects the most probable class for element e , as that yielding the highest value of density distribution functions at $h = h_e$. In other words, if $\rho_A(h_e) \geq \rho_B(h_e)$, element e is classified as A, otherwise it is classified as B.

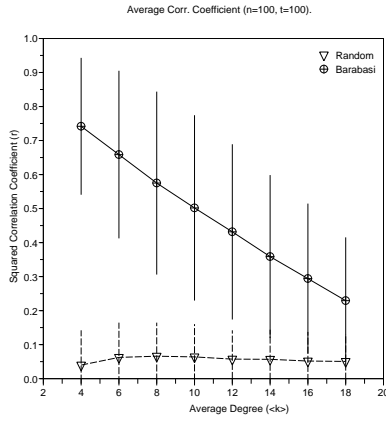


Figure 2: Average correlation coefficients for BA and ER models, taken from 2000 networks and considering every node of these networks.

The Density distribution functions are not always available, but can be estimated by using various methods. Here we consider non-parametric estimation from the respective normalized histogram obtained from the measurements. In order to improve the density distribution estimation, one can interpolate the histogram by using the Parzen windows method [12], which consists of convolving the histogram with a Gaussian distribution, resulting in an interpolated and smoother curve. Classification through Bayes map can then be obtained by considering the Pearson correlation coefficients obtained for the two studied models. An example of such histograms is depicted in figure 3. This figure includes the two histograms obtained for the BA and ER models for $\langle k \rangle = 4$ (a), 8 (b), 12 (c) and 16 (d). Note that the separation between the histograms decreases substantially with the increase of $\langle k \rangle$

4 Experimental Methodology

A software was developed in the Java language providing compatibility with any major operating system, specially in order to provide a graphical interface through which the subject can navigate along complex networks. A sequence of sets of networks with increasing average node degree were considered. Some parameters must be fixed prior to each navigation, including the model of the network (ER or BA), the number of nodes in each network, the average node degree of the initial set of networks², the average node degree for the last set and the number of networks per set. The navigation starts at a randomly selected

²A network set consists in a group of fixed average degree networks.

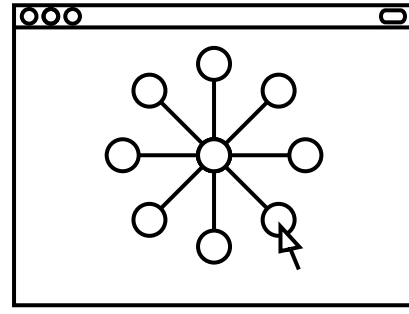


Figure 4: Navigation Screen: the subject navigates through the network while trying to determine whether it is random or scale free.

node, which becomes the current *central node*. The software records all the navigation actions taken by the subject for posterior analysis. Figure 4 shows a current node (center) and its immediate neighbors, where each node is represented by a circle and the edges are represented by lines linking these circles.

At each step, the subject is prompted to choose a node amongst the neighbors of the central node in order to continue the walk. The chosen node becomes the current central node, and the process is repeated. The subject can navigate until a specific number of steps is reached, prompting the user to make a decision about the model. No partial results of the experiment are presented to the subject during the navigation.

The subject can choose between either Barabási-Albert (BA) or Random model (ER). After the choice is made, it is stored and a new network is generated and showed to subject, repeating the process for every network of the set. After all networks in a set are navigated, the total number of correct choices is stored and the average degree is increased by one for the next set.

5 Virtual Agent Navigation

In order to have a comparison standard, and also to consider explicitly a model of navigation, an artificial agent has been developed which is capable to navigate through the same type of networks as humans. The adopted heuristics is described below and is schematized in Figure 5.

Two sets of 1000 networks (BA and ER) with node degree $N = 100$ and the same value of average degree are considered (figure 5:A). For every network, the agent is placed initially at a randomly chosen node and began its navigation while using an algorithm as described in [2], where the agent selects randomly a not yet visited neighbor node and makes it the new center. If only visited nodes exists in the neighbor, the agent use the same process to choose a new cen-

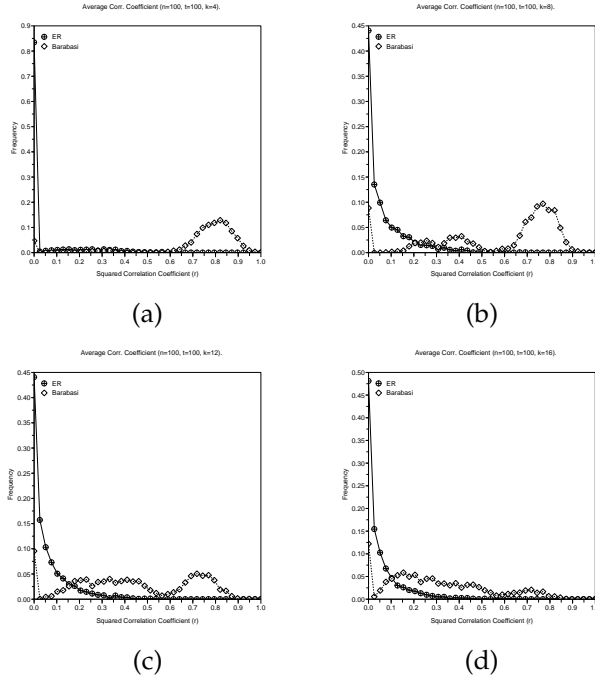


Figure 3: Average correlation coefficients for BA and ER models, considering every node in 2000 networks with $N=100$ and several values of $\langle k \rangle$

ter among those already visited nodes. The process is repeated until every network node has been visited.

For every new node, the number of already visited nodes, t , is stored and increased. A histogram is made (figure 5:B) considering the logarithmic of degree distribution along the network at each value of t , from which the respective Pearson correlation coefficient $c(t)$ is then estimated. The same process is repeated for every network of the same model type, obtaining a new set of histograms $H(t)$ from the distribution of correlation coefficient (figure 5:C). The same method is applied for the other model.

Because the correlation coefficients of BA networks tend to be higher than for ER, a meaningful decision regions can be created by considering the histograms of these measurements, indicating areas in a plane t versus $c(t)$ where the networks are more likely to correspond to BA or ER models (figure 5:D). Interestingly, different decision regions are obtained for sets of network with different values of average degree. A new set of networks is then navigated by the agent while considering the decision regions, obtaining the ratio of correct guesses for every new node visited (figure 5:F). Some of these maps, obtained by performing simulations on 2000 networks of each model are shown in figure 6.

The entire procedure is executed for networks with different values of average node degrees, resulting in a correct ratio surface on k vs. t plane.(figure 5:G). A simulation for 2000 networks resulted in correct ratio curves by the automated agents as shown in figure

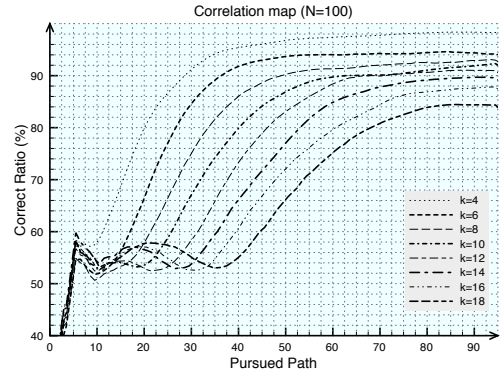


Figure 7: Correct ratio curves for the automated agents for 2000 simulated networks.

7 for several values of average degrees. As could be expected, the correct classification tends to improve along each curve for longer pursued paths. Also, It is clear from this figure that the correct classifications tend to progressively diminish for higher average degrees. These results shows that it becomes much more difficult to infer the nature of the network (i.e. ER or BA) when they are more dense (i.e. higher average node degree).

The results obtained for three subjects are shown in Figure 8. This figure shows the rate of correct classifications in terms of the average node degree. Recall that the classification was reached after a 30 random

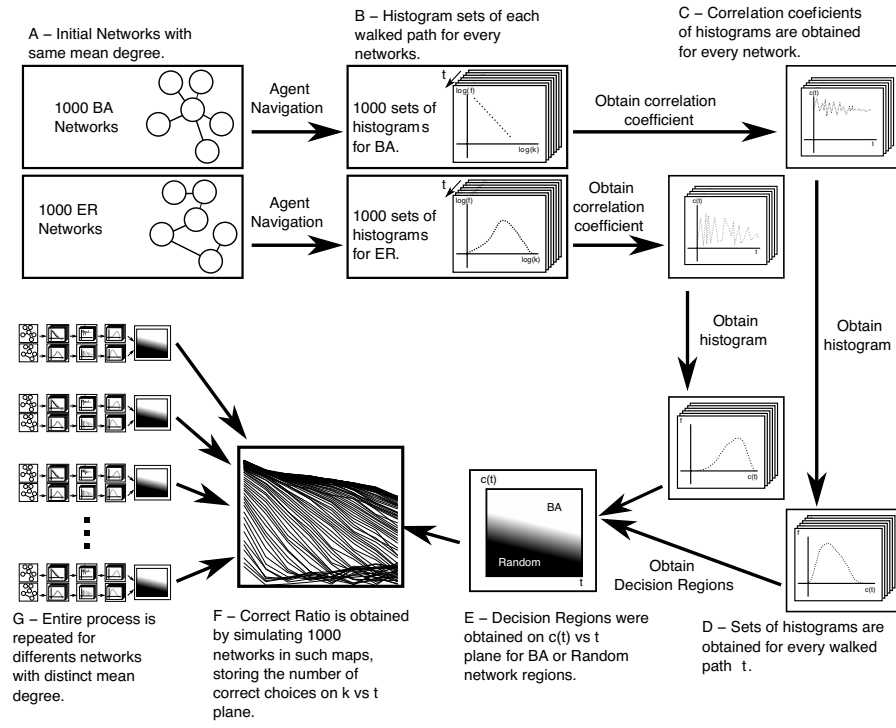


Figure 5: Schematic of main process to obtain decision regions and correct ratio for agent navigation.

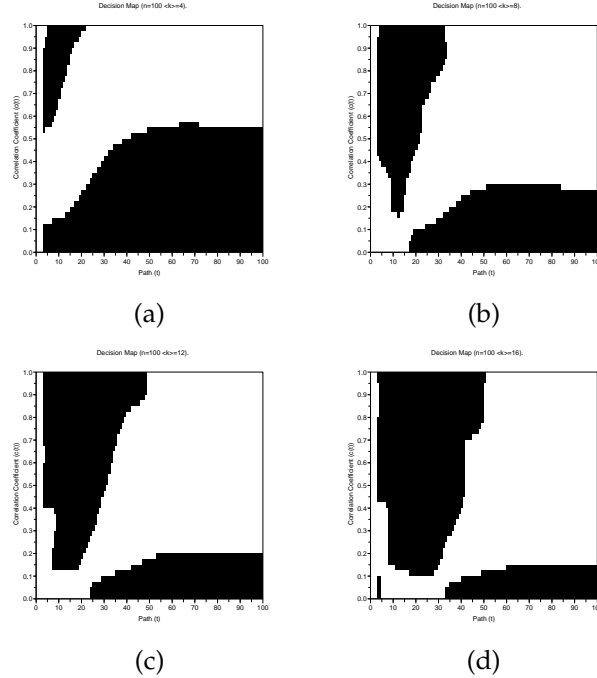


Figure 6: Decision regions obtained for the simulated models. The white region represents where BA model is most probable, while the black region does the same with respect to the ER model.

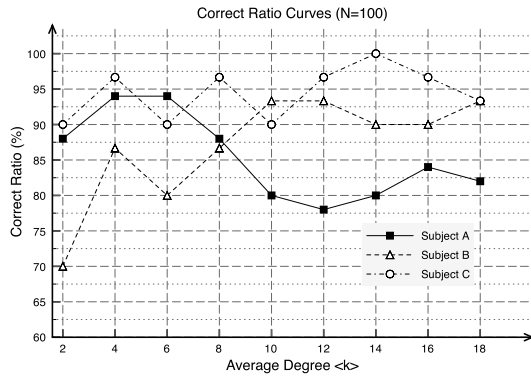


Figure 8: Several experiments obtained for three distinct subjects while considering a fixed number of walks of size 15.

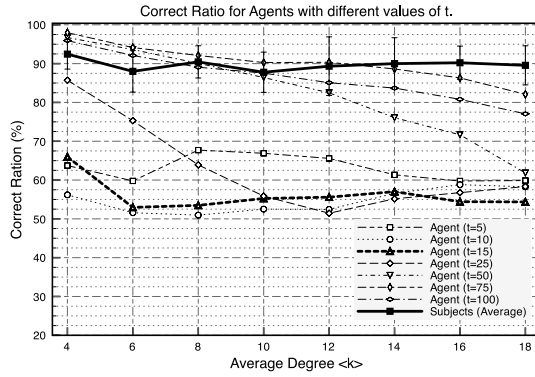


Figure 9: Autonomous agent correct ratio curves for several pursued paths (t) and average value and respective standard deviation for the results of experiments. (Comparable results are highlighted)

walks with 15 steps. Interestingly, unlike the automated case, the performance of the classification does not clearly diminish with the average degree. On the other hand, some fluctuations are observed for the correct classification ratio. The values of this ratio varied between 80 and 100%.

Figure 9 shows the average \pm standard deviations of the correct classification ratios in terms of the average node degree obtained for the automated and human agents. In the former case, seven curves are shown respective to different number of steps taken by the automated agents before making the decision. Note that a big change takes place for the automated agents when more than 40 steps are allowed before decision. This changes proceeds from less than 70% to over 85%. Figure 9 also shows the average \pm standard deviation obtained for the humans, whose average value fluctuates around 90%.

6 Concluding Remarks

This article has investigated the important problem of deciding on the type of network as one (automated agent or human subject) navigates along it. Two types of networks were considered: Barabási-Albert and Erdős-Rényi. The measurement of the network considered for the automated decision is the Pearson correlation coefficient extracted from the loglog distribution of node degrees. Bayesian decision theory was used in order to decide on the most likely type of network. A graphic-interactive interface was developed especially for human navigation, with the estimation of the network type being requested after 15 steps of the walk. The obtained results present a series of interesting features. First, both the automated agent and humans presented surprisingly good performance for identification of the type of network. Interestingly, such a performance tended to reduce in the case of the automated agent, while remaining constant (with some fluctuations) in the case of the humans. An abrupt change in the number of correct classification ratio was observed for the automated agents while moving from 40 to 70 steps.

All in all, the obtained results corroborate the ability of automated and human agents for discriminating between ER and BA complex networks with the same number of nodes and average degree. Additional investigations can be performed in order to identify which topological clues are being considered by the humans while trying to identify the type of the networks. In addition, it would be interesting to verify how the consideration of additional measurements of the networks (e.g. clustering coefficient, node correlations, shortest paths, etc.) may contribute for enhancing the performance of the autonomous agent.

Acknowledgment: Luciano da F. Costa thanks CNPq for partial sponsorship.

References

- [1] L. da F. Costa, physics/0601118, 2006
- [2] L. da F. Costa, physics/0604193, 2006
- [3] B. Tadic, Eur. Phys. J. B, 23 221 2001
- [4] B. Tadic, cond-mat/0310014, 2003
- [5] E. M. Bollt and D. ben Avraham, cond-mat/0409465, 2004
- [6] J. D. Noh and H. Rieger, cond-mat/0310344, 2004
- [7] M.E.J. Newman, SIAM Review, 45 167256 2003
- [8] Albert and A.-L. Barabási, Rev. Mod. Phys., 45 47-97 2002

- [9] L. da F. Costa, F. A. Rodrigues, G. Travieso e P. Villas Boas, Characterization of complex networks: A survey of measurements cond-mat/0505185, 2005
- [10] Cohen J, Statistical Power Analysis for the Behavioral Sciences (Lawrence Erlbaum Associates) 1988
- [11] L. da F. Costa, R. M. Cesar Jr. , Pattern Classification and Scene Analysis (CRC; 1 edition) 2000
- [12] Duda, R. O. and Hart, P. E. , Shape Analysis and Classification: Theory and Practice (Wiley) 1973